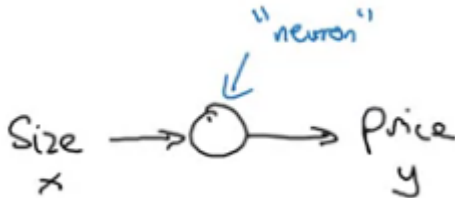


Week 1

What is a neural network?

- Neuron = takes an input and transforms it into output through a linear function



- Neural network (NN) = stacking neurons
 - one neuron for each recognized/classified feature
- Every input layer is connected to all neurons in a hidden layer

Supervised learning

- CNN = images
- RNN (recurrent) = temporal series (sequence data)
- Structured data = can be shown in a table
 - Every feature has a well-defined meaning
- Unstructured data (audio, image, text)

Why is deep learning taking off?

- m = size of the training set
- smaller training set = order of NN algorithms are not well-defined
 - we need good skills to determine features manually
- Bigger dataset = larger NN are better

Week 2

Binary Classification

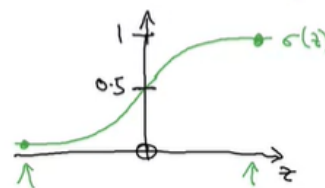
- I want to iterate all datasets without a for loop
 - forward and backward propagation
- (x, y) pair > x is a value, and y is the label
- the dataset has 2 lines (x and y) and m columns

Logistic regression

Given x , want $\hat{y} = P(y=1|x)$
 $x \in \mathbb{R}^{n_x}$ $0 \leq \hat{y} \leq 1$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

Output $\hat{y} = \sigma(\underbrace{w^T x + b}_z)$



- The sigmoid function goes from 0 to 1 = adequate to indicate a probability
- Objective = choose ω and b such that the sigma returns the probability of x return y as label

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

If z large $\sigma(z) \approx \frac{1}{1+0} = 1$

If z large negative number

$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx \frac{1}{1 + \text{Big num}} \approx 0$$

Andrew Ng

Logistic regression cost functions

$$L(y, \hat{y}) = - (y \log \hat{y} + (1-y) \log (1-\hat{y})) \leftarrow$$

If $y=1$: $L(y, \hat{y}) = -\log \hat{y} \leftarrow$ want $\log \hat{y}$ large, want \hat{y} large.
 If $y=0$: $L(y, \hat{y}) = -\log (1-\hat{y}) \leftarrow$ want $\log (1-\hat{y})$ large ... want \hat{y} small

Cost function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$

Gradient descent

- dJ/dw = gradient

Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

}

$w := w - \alpha dw$

learning rate

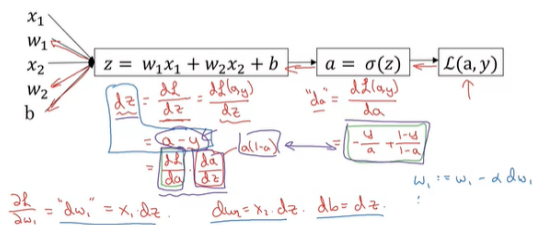
- $J(w, b)$ = dependent variable

Computation graph

- J is the final step of calculation (cost function)
- We have w and b and compute J (forward propagation)
- We then have J and compute the derivatives in relation to w and b (backward propagation)

Logistic regression gradient descent

Logistic regression derivatives



Gradient descent on m Examples

Logistic regression on m examples

$$J=0; dw_1=0; dw_2=0; db=0$$

For $i=1$ to m

$$z^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$\frac{dz^{(i)}}{dw_1} = a^{(i)} - y^{(i)}$$

$$\frac{dz^{(i)}}{dw_2} = a^{(i)} - y^{(i)}$$

$$\frac{dz^{(i)}}{db} = 1$$

$$dw_1 += x_1^{(i)} \frac{dz^{(i)}}{dw_1}$$

$$dw_2 += x_2^{(i)} \frac{dz^{(i)}}{dw_2}$$

$$db += \frac{dz^{(i)}}{db}$$

$J/=m$; $dw_1/=m$; $dw_2/=m$; $db/=m$

$dw_1 = \frac{\partial J}{\partial w_1}$

$w_1 := w_1 - \alpha \frac{\partial J}{\partial w_1}$

$w_2 := w_2 - \alpha \frac{\partial J}{\partial w_2}$

$b := b - \alpha \frac{\partial J}{\partial b}$

Andrew Ng

- there are 2 for loops (m examples and n dimensions of x)

Vectorization

- saves time to compute $z = w^T x + b$

Logistic regression derivatives

$$J = 0; dw_1 = 0; dw_2 = 0; db = 0$$

for $i = 1$ to m :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$\frac{dz^{(i)}}{dw} = x^{(i)} (a^{(i)} - y^{(i)})$$

$$\frac{dz^{(i)}}{db} = 1$$

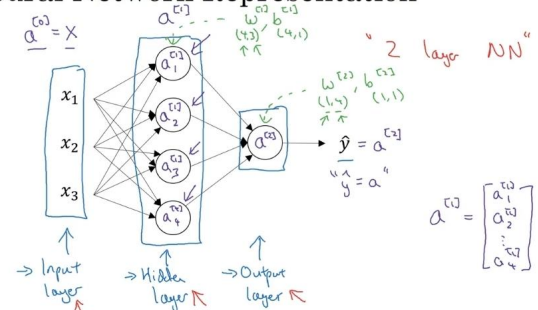
$$J = J/m; dw_1 = dw_1/m; dw_2 = dw_2/m; db = db/m$$

Week 3

Neural Network Representation

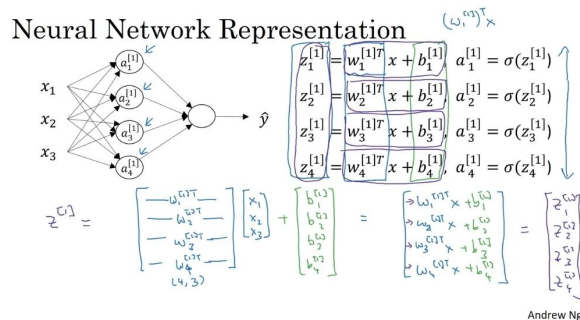
- Superscript with brackets = layer
- Superscript with parenthesis = element in the dataset
- $x = a[0]$ (activation parameters for the first layer)

Neural Network Representation



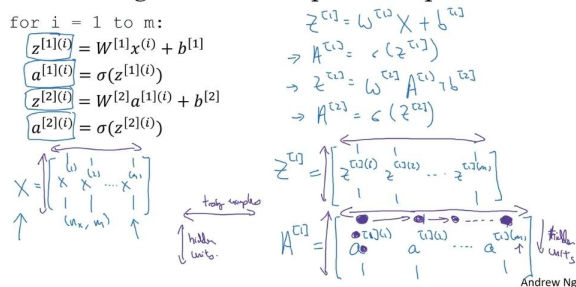
Computing a Neural Network's output

- Vectorizing logistic regression for multiple nodes (vertically stacked)
- Detail: this diagram shows for one input data x of 3 items each



Vectorizing across multiple examples

Vectorizing across multiple examples



Activation functions

- Sigmoid function may not be appropriate >> good for binary classification at the output layer (returns 0 or 1)
- tanh goes from -1 to 1 >> good for hidden layers
 - The mean of the values is zero >> centring data (see 2nd course)
- if z is very large or very small, tanh and sigmoid is approximately zero >> slows down learning
 - use ReLu function >> $dz = 1$ for z positive or $dz = 0$ for z negative
- Leaky Relu >> $dz = -k$, for k is small

Why use a non-linear activation function?

- If you use a linear function, the output of the NN is a linear function of the input X
 - So there's no difference in using one or more hidden layers
- But you can use it at the output layer if the predicted value is of the same nature as the input

Backpropagation with multiple layers and examples

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]}a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T}dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]}x^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]}A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$$

$$dZ^{[1]} = W^{[2]T}dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]}X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

Andrew Ng

Initiate with random values

- If all same units have the same function and the same initial weight, the result of the backdrop will be the same
- It's enough to initiate w non-zero and b thanks
- W must be small to avoid flat parts of sigmoid or tanh function >> could slow down learning

Week 4

Deep L-layer NN

- L = number of layers
- $n[L]$ = number of neurons in layer L
- First countable layer is not the input layer
 - input layer = 0
- $a[L]$ = activation variable of layer L

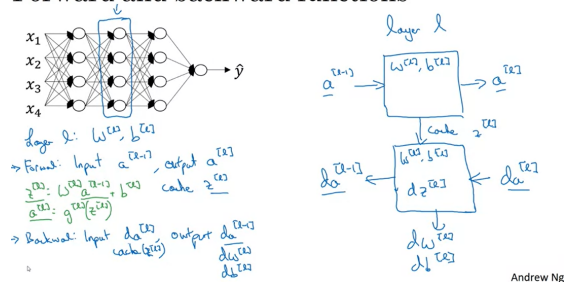
Forward propagation in a deep NN

$$\begin{aligned} z^{[l]} &= W^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

- Different examples = different columns

Building blocks in deep NN

Forward and backward functions



Hyperparameters

- Learning rate (alpha)
- How many iterations must I do
- How many hidden layers
- How many hidden units per layer
- What's the best activation functions